# Beyond Analysis

## Team Name: VoidMen

Abishek Mahesh  -  Prithvi Seshadri  -  Rishi Dinesh  -  Pran...

# Introduction

## Problem Statement

Head digital works has provided a dataset that falls in the domain of online skill-based gaming. This dataset contains a 22-dimensional feature space that describes the differ— aspect of customer behaviour. The goal is to predict the Y1 and Y2 values of a given customer, which represents the customer's current value and future value (temporal extrapolation) respectively

## Dataset

The dataset is divided into train and test sets. Each customer is identified by a UNIQUE_IDENTIFIER and a sequence of entries tagged by the SEQUENCE_NO. Various information attributing to the performance of the customer in the game has been prov
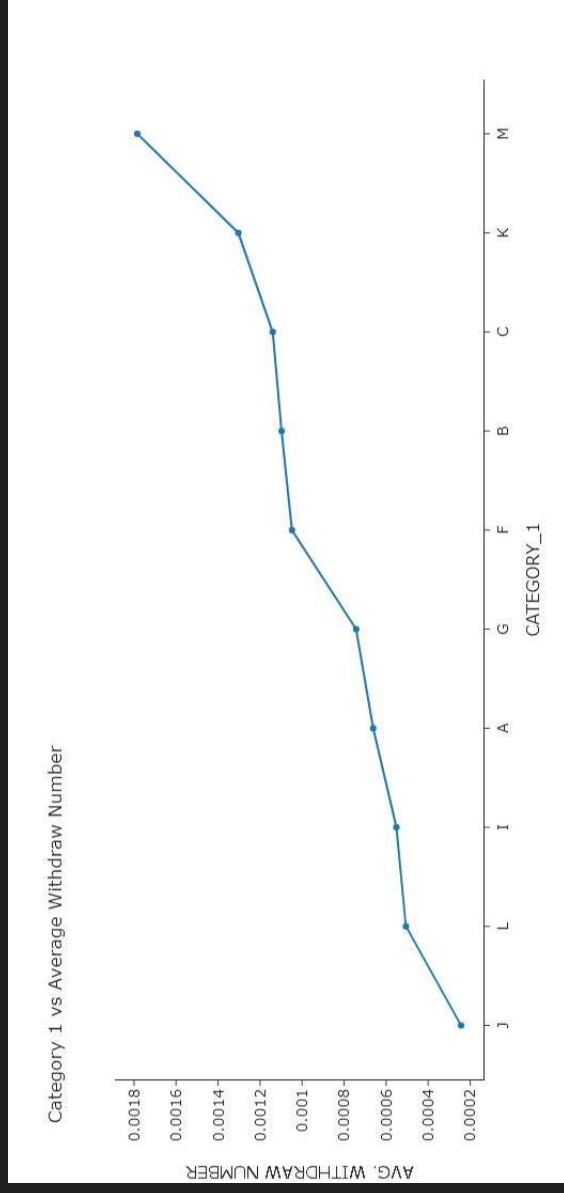
## Metrics

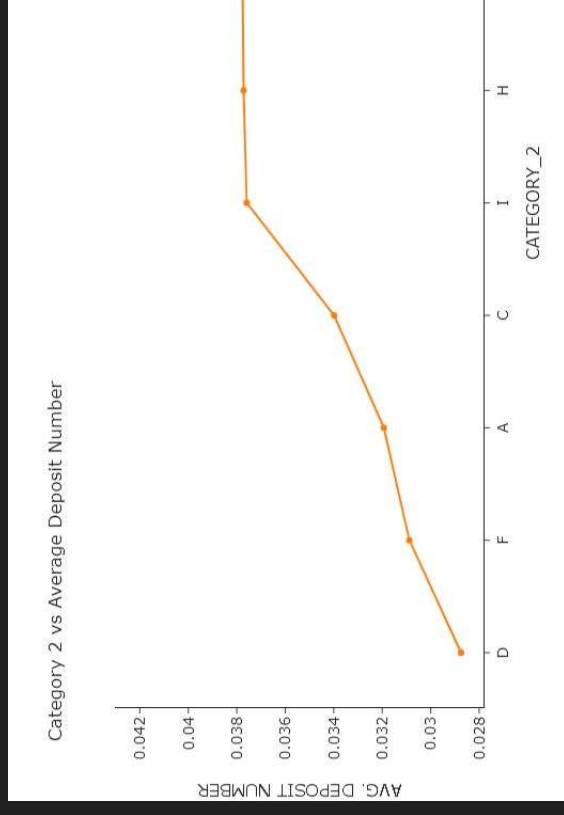Root Mean Squared Error(RMSE) is the primary metric used in our analysis.

# Exploratory Data Analysis



*Correlation Heatmap of Numerical Features*

# Exploratory Data Analysis

## CATEGORY_2 Visuali



Category 2 vs Average Deposit Number

## CATEGORY_1 Visualization



Category 1 vs Average Withdraw Number

# Exploratory Data Analysis

## Practice vs Live Data



## REVENUE Correlation Analysis

# Data Preparation

## Original

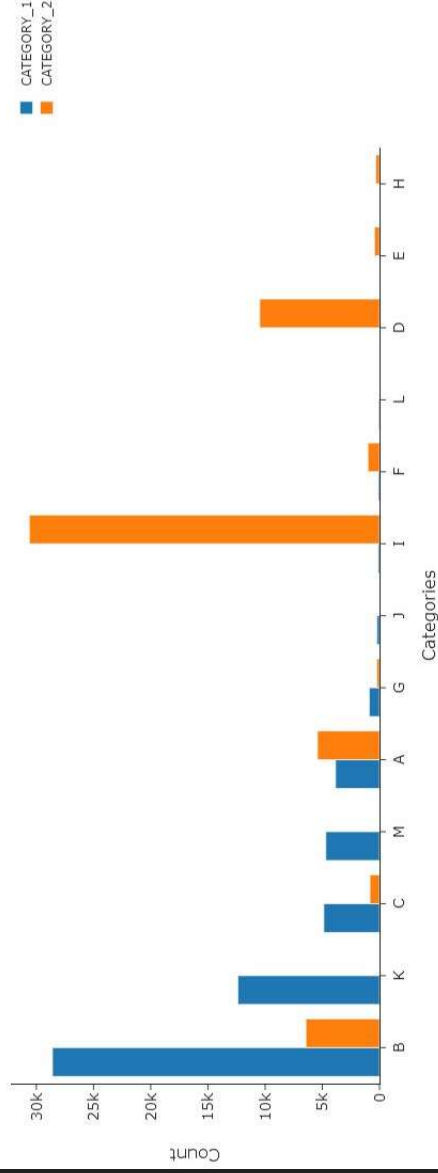| UNIQUE_IDENTIFIER | SEQUENCE_NO | STATUS_CHECK | CATEGORY_1 | CATEGORY_2 | ACTIVE_YN | ENTRY | REVENUE |
|---|---|---|---|---|---|---|---|
| 98481267304 | 1 | 0 | M | B | 1 | 0.000000 | 0.000000 |
| 98481267304 | 2 | 0 | M | B | 1 | 0.137350 | 0.011550 |
| 98481267304 | 3 | 0 | M | B | 1 | 0.158350 | 0.010425 |
| 98481267304 | 4 | 0 | M | .B | 1 | 0.444900 | 0.035850 |
| 98481267304 | 5 | 0 | M | B | 1 | 0.000000 | 0.000000 |
| 98481267304 | 6 | 0 | M | B | 1 | 0.000000 | 0.000000 |
| 98481267304 | 7 | 0 | M | B | 1 | 0.045050 | 0.002950 |

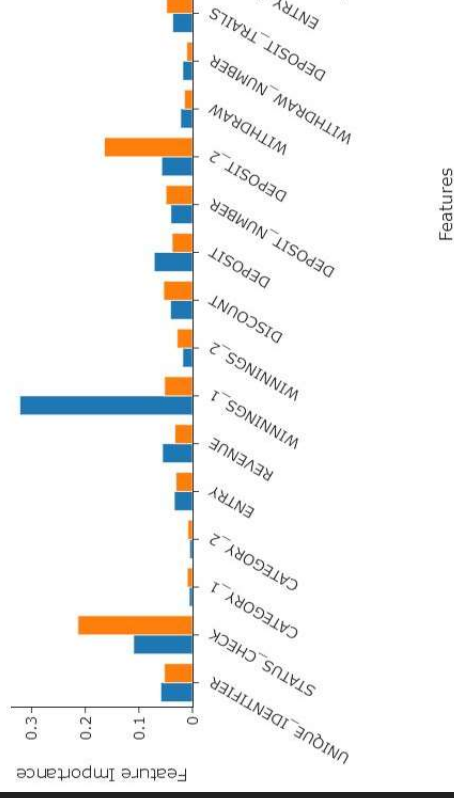| UNIQUE_IDENTIFIER | STATUS_CHECK | CATEGORY_1 | CATEGORY_2 | ENTRY | REVENUE | WINNINGS_1 | WINNINGS_2 | DISCOUNT | DEPOSIT | ... |
|---|---|---|---|---|---|---|---|---|---|---|
| 98481267304 | 0 | M | B | 0.098334 | 0.007531 | 0.043399 | 0.000000 | 0.000000 | 0.000714 | ... |
| 98481267698 | 1 | M | I | 31.392245 | 3.803991 | 25.940547 | 0.000000 | 0.866865 | 11.122807 | ... |
| 98481269325 | 0 | M | D | 0.018567 | 0.001624 | 0.010514 | 0.000000 | 0.005791 | 0.000278 | ... |
| 98481271512 | 0 | M | E | 0.747600 | 0.117320 | 0.025330 | 0.000000 | 0.240000 | 0.504000 | ... |
| 98481273023 | 0 | M | I | 0.500000 | 0.080000 | 0.000000 | 0.000000 | 0.000000 | 0.500000 | ... |

# Data Preprocessing

## Feature Importances (T...



Feature Importance Graph for Y1 and Y2

## Categories Frequency Chart



Distribution of Customer Categories

Methods Explored

# Final Ensemble Model

Voting Regressor(
Multilayer perceptron,
Random forest,
XGboost
)

- *An interactive visualization*

# Results

Train set evaluation Y1:

```
482/482 [==============================] - 0s 597us/step
MAE: 0.9461894717876128
MSE: 14.005275338171362
RMSE: 3.7423622670094323
R2 Square 0.8035067561731131
```

Train set evaluation Y2:

```
1926/1926 [==============================] -
MAE: 56.6194303064907
MSE: 9180.27927298143
RMSE: 95.8137739209838
R2 Square 0.7325900531845706
```

# Epilogue

## Conclusion

Various regression models such as RandomForest, Multilayer Perceptron and XGBoost were trained on a processed dataset. Results showed that an ensemble of the best performing models would yield the best score.

## Challenges

Understanding the dataset in the context of its domain was a hurdle. Tackling the temporal component of the dataset was also a difficulty.

## Future Scope

Using a deep learning approach for better feature selection and modelling